
A Probabilistic Representation for Deep Learning: Delving into The Information Bottleneck Principle

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The Information Bottleneck (IB) principle has recently attracted great attention to
2 explaining Deep Neural Networks (DNNs), and the key is to accurately estimate the
3 mutual information between a hidden layer and dataset. However, some unsettled
4 limitations weaken the validity of the IB explanation for DNNs. To address these
5 limitations and fully explain deep learning in an information theoretic fashion, we
6 propose a probabilistic representation for deep learning that allows the framework
7 to estimate the mutual information, more accurately than existing non-parametric
8 models, and also quantify how the components of a hidden layer affect the mutual
9 information. Leveraging the probabilistic representation, we take into account the
10 back-propagation training and derive two novel Markov chains to characterize the
11 information flow in DNNs. We show that different hidden layers achieve different
12 IB trade-offs depending on the architecture and the position of the layers in DNNs,
13 whereas a DNN satisfies the IB principle no matter the architecture of the DNN.

14 1 Introduction

15 Deep learning [18] has already achieved great success in numerous applications. Deep Neural
16 Networks (DNNs), however, are still commonly viewed as ‘black boxes’ [27]. Considerable efforts
17 have been devoted to explaining the internal mechanism of DNNs from various perspectives, such as
18 mathematics [5, 12], statistics [14, 20, 23], computer vision [37, 21], *etc.* Recently, the Information
19 Bottleneck (IB) principle has attracted attention in opening the ‘black boxes’ of DNNs [30, 33].

20 Given a joint distribution $P(X, Y)$, the IB principle posits a random variable $T = f(X)$ obeying the
21 Markov chain $Y \rightarrow X \rightarrow T$ and optimizes T by the IB Lagrangian [32, 31]

$$\min_{P(T|X)} I(X; T) - \beta I(Y; T), \quad (1)$$

22 where $f(\cdot)$ is an arbitrary function, $I(\cdot; \cdot)$ denotes mutual information, and the Lagrange multiplier
23 $\beta > 0$ controls the IB trade-off between compressing the input X and preserving the information
24 of the label Y . In the seminal work [30], Tishby *et al.* manifest the IB trade-off in every layer of
25 DNNs $= \{\mathbf{x}; \mathbf{t}_1; \dots; \mathbf{t}_L; \hat{\mathbf{y}}\}$ via studying $I(X; T_i)$ and $I(Y; T_i)$, where T_i is the random variable of
26 the i th hidden layer \mathbf{t}_i . Especially, the authors ascribe DNN generalization to the compression [29].

27 In the context of deterministic DNNs, recent works reveal some limitations of the IB principle for
28 explaining DNNs. Amjad *et al.* argue that the IB principle becomes an ill-posed optimization problem
29 due to $I(X; T_i) = \infty$ [1], and Kolchinsky *et al.* demonstrate that not every layer of DNNs satisfies a
30 strict IB trade-off, *i.e.*, different layers only differ in $I(X; T_i)$ but $I(Y; T_i)$ keeps consistent in all
31 layers [15]. In addition, Saxe *et al.* experimentally show that the compression does not occur in
32 DNNs with non-saturating activation functions, *e.g.*, the popular ReLU function [28], and Goldfeld
33 *et al.* doubt the causality between the generalization of DNNs and the compression [10, 7]. These
34 unsettled limitations greatly weakens the validity of the IB explanations for DNNs.

35 The key to examining the IB principle in DNNs is the accurate estimation of the mutual information.
 36 However, regarding DNNs as deterministic models hinders us from specifying the random variable
 37 T_i and the distribution $P(T_i)$, thus it is difficult to accurately estimate $I(X; T_i)$ and $I(Y; T_i)$. More
 38 specifically, in the absence of a clear definition of T_i , simply assuming the activations of t_i as the *i.i.d.*
 39 samples of T_i induces T_i being a continuous random variable and $I(X; T_i) = \infty$ in deterministic
 40 DNNs (see Appendix C in [28]). The complicated architecture of DNNs makes it challenging to
 41 specify $P(T_i)$. Therefore, most previous works have to indirectly estimate $P(T_i)$ via non-parametric
 42 models [35], such as the empirical distribution [30], Kernel Density Estimation (KDE) [28], and
 43 Gaussian convolution [10]. However, we experimentally confirm that classical non-parametric models
 44 derives poor mutual information estimation [24, 22] in DNNs, and one reason is because activations
 45 do not satisfy the *i.i.d.* prerequisite of non-parametric models (see Appendix G). In summary, the
 46 limitations mainly stem from the lack of an explicit probabilistic representation for deep learning.

47 The IB principle only formulates the information flow in DNNs = $\{\mathbf{x}, \mathbf{t}_1, \dots, \mathbf{t}_I, \hat{\mathbf{y}}\}$ after training,
 48 and the corresponding Markov chain (see Fig. 1 in [30])

$$Y \rightarrow X \rightarrow T_1 \cdots \rightarrow T_I \rightarrow \hat{Y} \quad (2)$$

49 indicates that the information of Y transfers to T_i in the forward direction and T_i receives the
 50 information of Y only via X . However, training DNNs by the back-propagation [25] implies that the
 51 information of Y transfers to T_i in the backward direction during training and retains information
 52 in T_i after training. Notably, Zhang *et al.* show that a DNN can fit labels well even using Gaussian
 53 noise as input to train the DNN [38], which implies that T_i can directly receive the information of Y .
 54 Hence, the IB principle does not comprehensively characterize the information flow in DNNs.

55 To address the above limitations and comprehensively explain DNNs in an information theoretic
 56 fashion, we introduce the probability space $(\Omega_{T_i}, \mathcal{F}, P_{T_i})$ [6] for the i th hidden layer t_i in DNNs.
 57 Compared to previous works, the probability space $(\Omega_{T_i}, \mathcal{F}, P_{T_i})$ enables us to: (i) accurately estimate
 58 $I(X; T_i)$ and $I(Y; T_i)$ via specifying T_i and $P(T_i)$, and (ii) quantify the effect of the architecture of
 59 t_i and the back-propagation on $I(X; T_i)$ and $I(Y; T_i)$ via explicitly modeling all the ingredients of t_i ,
 60 such as the activation function and the weights in a probabilistic way. To the best of our knowledge,
 61 this is the first time the probability space of a hidden layer in DNNs is as defined.

62 Leveraging $(\Omega_{T_i}, \mathcal{F}, P_{T_i})$, we derive information theoretic explanations for DNNs as follows:

- 63 • Two Markov chains¹ characterize the information flow in DNNs = $\{\mathbf{x}, \mathbf{t}_1, \dots, \mathbf{t}_I, \hat{\mathbf{y}}\}$

$$\begin{aligned} \bar{X} &\rightarrow T_1 \rightarrow \dots \rightarrow T_I \rightarrow \hat{Y} \\ T_1 &\leftarrow \dots \leftarrow T_I \leftarrow \hat{Y} \leftarrow Y. \end{aligned} \quad (3)$$

- 64 • Different hidden layers manifest different IB trade-offs depending on the architecture and
 65 the position of hidden layers in DNNs.
- 66 • A DNN satisfies the IB principle no matter the architecture of the DNN.

67 **Preliminaries.** $P(X, Y) = P(X)P(Y|X)$ is an unknown joint distribution between X and Y . A
 68 dataset $\mathcal{D} = \{(\mathbf{x}^j, y^j) | \mathbf{x}^j \in \mathbb{R}^M, y^j \in \mathbb{Z}\}_{j=1}^J$ consists of J *i.i.d.* samples generated from $P(X, Y)$
 69 with finite L labels, i.e., $y^j \in \{1, \dots, L\}$. In the context of supervised learning, we focus on
 70 feedforward fully connected DNNs = $\{\mathbf{x}, \mathbf{t}_1, \dots, \mathbf{t}_I, \hat{\mathbf{y}}\}$, i.e., Multi-Layer Perceptions (MLPs) [8]
 71 for the image classification task. Without loss of generality, we use the MLP = $\{\mathbf{x}, \mathbf{t}_1, \mathbf{t}_2, \hat{\mathbf{y}}\}$ with
 72 the cross-entropy loss ℓ_{CE} for most theoretical derivations. In addition, $H(\cdot)$ denotes entropy.

73 In the MLP, \mathbf{t}_1 and \mathbf{t}_2 have N and K neurons, respectively, and $\mathbf{t}_1 = \{t_{1n} = \sigma_1[\langle \boldsymbol{\omega}_n^{(1)}, \mathbf{x} \rangle]\}_{n=1}^N$,
 74 where $\langle \boldsymbol{\omega}_n^{(1)}, \mathbf{x} \rangle = \sum_{m=1}^M \omega_{mn}^{(1)} \cdot x_m + b_{1n}$ is the n th dot-product given the weight $\omega_{mn}^{(1)}$ and the bias
 75 b_{1n} , and $\sigma_1(\cdot)$ denotes an activation function, e.g., ReLU. Similarly, $\mathbf{t}_2 = \{t_{2k} = \sigma_2[\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle]\}_{k=1}^K$,
 76 where $\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle = \sum_{n=1}^N \omega_{nk}^{(2)} \cdot t_{1n} + b_{2k}$. The output layer $\hat{\mathbf{y}}$ is softmax with L nodes

$$\hat{\mathbf{y}} = \{\hat{y}_l = \frac{1}{Z_Y} \exp[\langle \boldsymbol{\omega}_l^{(3)}, \mathbf{t}_2 \rangle] = \frac{1}{Z_Y} \exp[g_l(\mathbf{t}_2(\mathbf{t}_1(\mathbf{x})))]\}_{l=1}^L, \quad (4)$$

77 where $\langle \boldsymbol{\omega}_l^{(3)}, \mathbf{t}_2 \rangle = \sum_{k=1}^K \omega_{kl}^{(3)} \cdot t_{2k} + b_{yl}$ and $Z_Y = \sum_{l=1}^L \exp[\langle \boldsymbol{\omega}_l^{(3)}, \mathbf{t}_2 \rangle]$ is the partition function.

¹In which the virtual random variable \bar{X} has all the information of X except Y , namely $H(\bar{X}) = H(X|Y)$.

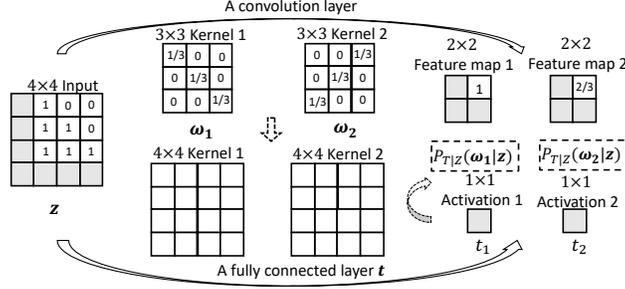


Figure 1: Given a 4×4 input z , a fully connected layer t is equivalent to a convolution layer with 4×4 convolution kernels. The definition of convolution (Chapter 9.1 in [11]) implies that the 4×4 weights ω_1 and ω_2 define two global features, and the two activations t_1, t_2 indicate the cross-correlation between ω_1, ω_2 and z , respectively. $P_{T|Z}(\omega_1|z)$ and $P_{T|Z}(\omega_2|z)$ measure the probability of ω_1 and ω_2 being recognized as the feature with the largest cross-correlation to z , respectively.

78 2 A probabilistic representation for deep learning

79 To accurately estimate $I(X; T_i)$ and $I(Y; T_i)$, in this section, we specify the probability space [6] for
80 a fully connected layer and derive the probabilistic explanations of the entire MLP.

81 It is known that a convolution kernel (namely the weights of convolution) defines a local feature,
82 and a convolution operation derives a feature map to measure the cross-correlation between the
83 local feature and input in a receptive field (Chapter 9.1 in [11]). Notably, a fully connected layer
84 is equivalent to a convolution layer with the kernel size having the same dimension as input. Thus
85 the weights of a neuron can be viewed as a global feature, and a fully connected layer with multiple
86 neurons derives activations to measure the cross-correlation between the multiple global features and
87 the input. The cross-correlation explanation for a fully connected layer is visualized in Figure 1.

88 Assuming that a fully connected layer t consists of N neurons $\{t_n = \sigma[\langle \omega_n, z \rangle]\}_{n=1}^N$, where $z \in \mathbb{R}^M$
89 is the input of t , $\langle \omega_n, z \rangle = \sum_{m=1}^M \omega_{mn} \cdot z_m + b_n$ is the dot-product between z and ω_n , and $\sigma(\cdot)$
90 is an activation function. Based on the cross-correlation explanation, the behavior of t is to measure
91 the cross-correlations between z and the N possible features defined by the the weights $\{\omega_n\}_{n=1}^N$.
92 In the context of pattern recognition [34], we define a virtual random process or ‘experiment’ as t
93 recognizing one of the patterns/features with the largest cross-correlation to z from the N possible
94 features. The experiment characterizes the behavior of t (*i.e.*, before recognizing the features with
95 the largest cross-correlation, t must measure the cross-correlations between z and all the N possible
96 features) while meets the requirement of the ‘experiment’ definition (*i.e.*, only one outcome will
97 occur on each trial of the experiment [6]). The probability space $(\Omega_T, \mathcal{F}, P_T)$ is defined as follows:

98 **Definition 1.** $(\Omega_T, \mathcal{F}, P_T)$ consists of three components: the sample space Ω_T has N possible
99 outcomes (features) $\{\omega_n = \{\omega_{mn}\}_{m=1}^M\}_{n=1}^N$ defined by the weights² of the N neurons; the event
100 space \mathcal{F} is the σ -algebra; and the probability measure P_T is a Gibbs distribution [19] to quantify the
101 probability of ω_n being recognized as the feature with the largest cross-correlation to z .

102 Taking into account the randomness of z , the conditional distribution $P_{T|Z}$ is formulated as

$$P_{T|Z}(\omega_n|z) = \frac{1}{Z_T} \exp(t_n) = \frac{1}{Z_T} \exp[\sigma(\langle \omega_n, z \rangle)], \quad (5)$$

103 where Z is the random variable of z and $Z_T = \sum_{n=1}^N \exp(f_n)$ is the partition function.

104 $(\Omega_T, \mathcal{F}, P_T)$ clearly explains all the ingredients of t in a probabilistic fashion. The n th neuron
105 defines a global feature by the weights w_n and the activation $t_n = \sigma(\langle \omega_n, z \rangle)$ measures the cross-
106 correlation between w_n and z . The Gibbs distribution $P_{T|Z}$ indicates that if w_n has the higher
107 activation, *i.e.*, the larger cross-correlation to z , it has the larger probability being recognized as
108 the feature with largest cross-correlation to z . For instance, if $z \in \mathbb{R}^{16}$ and t includes $N = 2$
109 neurons, then $\Omega_T = \{\omega_1, \omega_2\}$ defines two possible outcomes (features), where $\omega_n = \{\omega_{mn}\}_{m=1}^{16}$.
110 $\mathcal{F} = \{\emptyset, \{\omega_1\}, \{\omega_2\}, \{\omega_1, \omega_2\}\}$ means that neither, one, or both of the features are recognized
111 by t given z , respectively. $P_{T|Z}(\omega_1|z)$ and $P_{T|Z}(\omega_2|z)$ are the probability of ω_1 and ω_2 being
112 recognized as the feature with the largest cross-correlation to z , respectively.

²We do not take into account the scalar value b_n for defining Ω_T , as it not affects the feature defined by ω_n .

113 $(\Omega_T, \mathcal{F}, P_T)$ explains the representation ability of deep learning. Compared to Restricted Boltzmann
 114 Machines (RBMs) [26] simply using binary units to indicate features being recognized or not given
 115 input, the Gibbs distribution³ $P_{T|Z}(\omega_n|z)$ measures the probability of ω_n being recognized with
 116 the largest cross-correlation to z , i.e., it characterizes the relation between features and input more
 117 accurately. Moreover, Equation 5 shows that $t_n = \sigma(\langle \omega_n, z \rangle)$ is the negative energy function [19] of
 118 the Gibbs distribution, thus $P_{T|Z}(\omega_n|z)$ can be derived as long as $\sigma(\langle \omega_n, z \rangle)$ are known because
 119 the energy function is the sufficient statistics [2] of the Gibbs distribution. That enables subsequent
 120 hidden layers to generate high-level features of input via directly processing the activations $\{t_n\}_{n=1}^N$,
 121 thus deep learning can form a hierarchical structure to represent much complex features.

122 $(\Omega_T, \mathcal{F}, P_T)$ answers a fundamental question: which component of a hidden layer contains the
 123 information of the layer? Since ω_n defines Ω_T , the weights contain all the information of a layer. In
 124 particular, since the activation $t_n = \sigma(\langle \omega_n, z \rangle)$ is a function of ω_n , the data processing inequality
 125 [4] indicates that the information of t_n is no more than the information of ω_n . Simulations in Section
 126 4.2 demonstrate that if activations do not correctly characterize the cross-correlation between weights
 127 and input, activations contain less information than weights do.

128 Based on $(\Omega_T, \mathcal{F}, P_T)$, we define the random variable T as follows:

129 **Definition 2.** Given the fully connected layer \mathbf{t} , we define the random variable $T : \Omega_T \rightarrow E_T$ as

$$T(\omega_n) \triangleq n, \quad (6)$$

130 where the measurable space $E_T = \{1, \dots, N\}$.

131 Since Ω_T is composed of finite N possible outcomes, T is a discrete random variable. Notably, the
 132 one-to-one correspondence between ω_n and n indicates

$$P_{T|Z}(\omega_n|z) = P_{T|Z}(n|z). \quad (7)$$

133 If not considering the back-propagation training, the weights (namely Ω_{T_i}) of each layer are fixed.
 134 Thus T_{i+1} entirely depends on T_i and the MLP = $\{\mathbf{x}; \mathbf{t}_1; \mathbf{t}_2; \hat{\mathbf{y}}\}$ forms a Markov chain

$$X \rightarrow T_1 \rightarrow T_2 \rightarrow \hat{Y}. \quad (8)$$

135 Based on the corresponding joint distribution $P(\hat{Y}, T_2, T_1|X) = P(T_1|X)P(T_2|T_1)P(\hat{Y}|T_2)$ and
 136 Definition 2, we derive a probabilistic explanation for the entire MLP, which is summarized in
 137 Theorem 1. The detailed derivation is presented in Appendix B.

138 **Theorem 1.** The MLP = $\{\mathbf{x}; \mathbf{t}_1; \mathbf{t}_2; \hat{\mathbf{y}}\}$ formulates a conditional Gibbs distribution

$$P_{\hat{Y}|X}(l|\mathbf{x}) = \sum_{k=1}^K \sum_{n=1}^N P(\hat{Y} = l, T_2 = k, T_1 = n|X = \mathbf{x}) = \frac{1}{Z_{\text{MLP}}(\mathbf{x})} \exp[g_l(\mathbf{t}_2(\mathbf{t}_1(\mathbf{x})))] \quad (9)$$

139 where $Z_{\text{MLP}}(\mathbf{x}) = \sum_{l=1}^L \sum_{k=1}^K \sum_{n=1}^N P_{\hat{Y}, T_2, T_1|X}(l, k, n|\mathbf{x})$ is the partition function.

140 Since $P_{\hat{Y}|X}(l|\mathbf{x})$ exactly equals the output \hat{y}_l of the MLP, namely Equation (4), we conclude that
 141 the entire architecture of the MLP forms a family of Gibbs distribution $P_{\hat{Y}|X}(l|\mathbf{x})$. In general, the
 142 back-propagation updates a weight ω based on the gradient of ℓ_{CE} with respect to ω ,

$$\omega(s+1) = \omega(s) - \alpha \cdot \frac{\partial \ell_{\text{CE}}}{\partial \omega(s)} = \omega(s) - \alpha \cdot \frac{\partial \text{KL}[P(Y|X)||P(\hat{Y}|X)]}{\partial \omega(s)}, \quad (10)$$

143 where s is the index of training iteration, α is the training rate, and $\text{KL}[\cdot||\cdot]$ is the KL-divergence.

144 Figure 2 summarizes the probabilistic explanation for deep learning based on the MLP. In general,
 145 a single learning iteration, an epoch, consists of two phases: training and inference (after training).
 146 During inference, the MLP bridges X and \hat{Y} via multiple intermediate features Ω_{T_1} , Ω_{T_2} , and $\Omega_{\hat{Y}}$
 147 defined by weights, and formulates the statistical relation between \hat{Y} and X as a family of conditional
 148 Gibbs distribution $P(\hat{Y}|X)$. During training, the back-propagation updates weights to learn optimal
 149 intermediate features for searching an optimal $P(\hat{Y}|X)$ to accurately approximate $P(Y|X)$.

³Recent works about Gibbs explanations for a hidden layer are discussed in Appendix A.

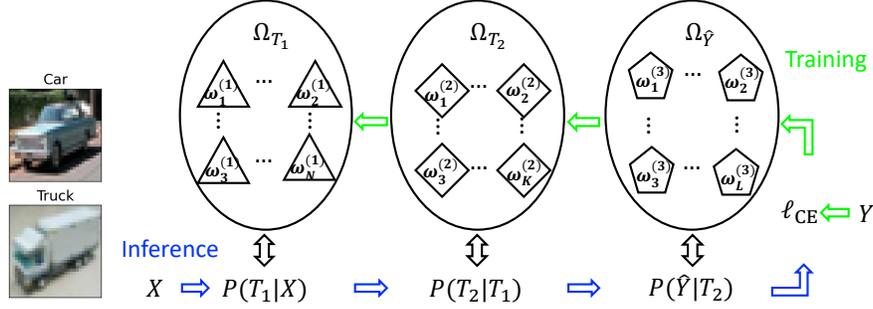


Figure 2: The visualization of the probabilistic explanation for deep learning based on the MLP.

150 3 The information theoretic explanations for deep learning

151 To address the limitations of existing IB explanations, this section proposes some novel information
 152 theoretic explanations for DNNs based on the proposed probabilistic representation.

153 **Proposition 1.** The mutual information between a fully connected layer and dataset is finite.

$$I(X; T) < \infty. \quad (11)$$

154 *Proof:* Definition 2 shows $E_T = \{1, \dots, N\}$. Thus T is a discrete random variable and $H(T) < \infty$,
 155 thereby $I(X; T) \leq H(T) < \infty$.

156 Proposition 1 circumvents the infinite mutual information problem. In the absence of a clear definition
 157 $T : \Omega_T \rightarrow E_T$, most previous works [28, 3, 1] simply viewing the activation t_n as the sample of T ,
 158 namely $t_n \in E_T = \mathbb{R}$, implies T being continuous and gives rise to the infinite mutual information
 159 problem in deterministic DNNs. However, $(\Omega_T, \mathcal{F}, P_T)$ indicates that t_n actually is a variable
 160 measuring the cross-correlation between w_n and z rather than the sample of T , namely $t_n \notin E_T$.

161 **Theorem 2.** The information of Y flows into the MLP in the backward direction during training

$$T_1 \leftarrow T_2 \leftarrow \hat{Y} \leftarrow Y. \quad (12)$$

162 *Proof:* First, since Ω_T is defined by ω in $(\Omega_T, \mathcal{F}, P_T)$ and Equation (10) shows that $\omega(s+1)$ is
 163 determined by all the previous gradients $\{\frac{\partial \ell_{CE}}{\partial \omega^{(s)}}\}_{s=1}^S$, and $\omega(0)$ is randomly initialized and α is a
 164 constant, we can derive that Ω_T is determined by $\frac{\partial \ell_{CE}}{\partial \omega}$. Second, based on the back-propagation, the
 165 relation between gradients in two adjacent layers in the MLP $= \{\mathbf{x}; \mathbf{t}_1; \mathbf{t}_2; \hat{\mathbf{y}}\}$ is formulated as

$$\begin{aligned} \frac{\partial \ell_{CE}^*}{\partial \omega_{kl}^{(3)}} &= [P_{\hat{Y}|X}(l|\mathbf{x}) - P_{Y|X}(l|\mathbf{x})] \cdot t_{2k}, \\ \frac{\partial \ell_{CE}^\circ}{\partial \omega_{nk}^{(2)}} &= \sum_{l=1}^L \frac{\partial \ell_{CE}^*}{\partial \omega_{kl}^{(3)}} \cdot \omega_{kl}^{(3)} \cdot \frac{\sigma_2'(\langle \omega_k^{(2)}, \mathbf{t}_1 \rangle)}{f_{2k}} \cdot t_{1n}, \quad \frac{\partial \ell_{CE}^\circ}{\partial \omega_{mn}^{(1)}} = \sum_{k=1}^K \frac{\partial \ell_{CE}^\circ}{\partial \omega_{nk}^{(2)}} \cdot \omega_{nk}^{(2)} \cdot \frac{\sigma_1'(\langle \omega_n^{(1)}, \mathbf{x} \rangle)}{t_{1n}} \cdot x_m. \end{aligned} \quad (13)$$

166 Equation 13 shows that $\frac{\partial \ell_{CE}^*}{\partial \omega^{(3)}}$ is a function of $P_{Y|X}(l|\mathbf{x})$ and $\frac{\partial \ell_{CE}^\circ}{\partial \omega^{(i)}}$ is a function of $\frac{\partial \ell_{CE}^*}{\partial \omega^{(i+1)}}$, where
 167 $\omega^{(3)}$ denotes the weight of $\hat{\mathbf{y}}$. The two points above enable us to derive that Ω_{T_i} is a function of $\Omega_{T_{i+1}}$
 168 and $\Omega_{\hat{Y}}$ is a function of $P(Y|X)$. Based on Definition 2, we can further derive that T_i is a function
 169 of T_{i+1} and \hat{Y} is a function of Y , i.e., $T_1 \leftarrow T_2 \leftarrow \hat{Y} \leftarrow Y$. (See the detailed proof in Appendix C).

170 Theorem 2 is consistent with the prevailing explanation for deep learning. LeCun *et al.* show that
 171 deep learning exploits the hierarchical property of signals [18], i.e., the layers farther from output
 172 learn lower-level features, such as edges, whereas the layers closer to output assemble lower-level
 173 features into the higher-level features corresponding to labels (see Figure 2 in [37]). Notably, since
 174 lower-level features commonly exist in signals with different labels (e.g., lower-level features, such
 175 as the edges of the vehicle frame and the circular contour of wheels, exist in both the car and the
 176 truck classes in the CIFAR-10 dataset [16] in Figure 2), lower-level features do not contain much
 177 information of labels. Therefore, the layers farther from output do not have much information of
 178 labels, which is consistent with the Markov chain $T_1 \leftarrow T_2 \leftarrow \hat{Y} \leftarrow Y$.

179 Since all the information of Y stems from X (i.e., $H(Y) = I(X; Y)$ proven in Appendix D),
 180 Theorem 2 implies that partial information of X flows into the MLP in the backward direction during
 181 training. Equation (2) shows the information of X flowing into the MLP in the forward direction
 182 during inference. Overall, the information of X flows in the backward and forward directions during
 183 training and inference, respectively. As a result, the Markov chain, Equation (2), proposed by recent
 184 works could not fully characterize the information flow of X in the MLP in each epoch. In other
 185 words, $I(X; T_i)$ is not necessarily greater than $I(X; T_{i+1})$ in the MLP in each epoch.

186 Equation (2) shows that T_i receives the information of Y via X during inference. Theorem 2 shows
 187 that T_i also directly receives information of Y during training, because the back-propagation updates
 188 weights (i.e., Ω_{T_i}) based on the label Y . Thus Equation (2) cannot fully characterize the information
 189 flow of Y in the MLP in each epoch, when we take into account the back-propagation training.

190 To fully characterize the information flow in the MLP in each epoch, we introduce Corollary 1.

191 **Corollary 1.** The information flow in the MLP can be characterized by two Markov chains as

$$\begin{aligned} \bar{X} &\rightarrow T_1 \rightarrow T_2 \rightarrow \hat{Y} \\ T_1 &\leftarrow T_2 \leftarrow \hat{Y} \leftarrow Y. \end{aligned} \quad (14)$$

192 The virtual random variable \bar{X} contains all the information of X except Y , i.e., $H(\bar{X}) = H(X|Y)$.

193 *Proof of the first Markov chain:* Since \bar{X} does not have any information of Y , it can only flow into
 194 the MLP in the forward direction during inference. Again since \bar{X} does not have any information of
 195 Y , the information flow of Y during training will not affect the information flow of \bar{X} . Therefore,
 196 $\bar{X} \rightarrow T_1 \rightarrow T_2 \rightarrow \hat{Y}$ characterizes the information flow of \bar{X} in both training and inference phases.

197 *Proof of the second Markov chain:* Since the weights are fixed after training, the sample space and
 198 the distribution of hidden layers are fixed after training. Therefore, the information of Y transferred
 199 into hidden layers during training will retain there after training (i.e., during inference). In addition,
 200 Definition 1 indicates that a fully connected layer $\mathbf{t} = \{t_n = \sigma(\langle \boldsymbol{\omega}_n, \mathbf{z} \rangle)\}_{n=1}^N$ measures the cross-
 201 correlation between $\boldsymbol{\omega}_n^{(1)}$ and \mathbf{z} during inference, thus $\{\boldsymbol{\omega}_n^{(1)}\}_{n=1}^N$ can be viewed as a representation
 202 of Z . As a result, even though Z has all the information of Y , the information of Y that \mathbf{t} can learn
 203 from Z is determined by how much information of Y the representation $\{\boldsymbol{\omega}_n^{(1)}\}_{n=1}^N$ has. Overall, the
 204 information flow of Y during inference will be the same as that during training. Based on Theorem 2,
 205 we conclude that $T_1 \leftarrow T_2 \leftarrow \hat{Y} \leftarrow Y$ characterizes the information flow of Y in the MLP in both
 206 training and inference phases. Detailed derivations and explanations are presented in Appendix E.

207 To quantify how much information of X and Y is learned by the MLP, we introduce Corollary 2.

208 **Corollary 2.** The mutual information between dataset and the entire MLP can be expressed as

$$\begin{aligned} I(X; T_{\text{MLP}}) &= I(\bar{X}; T_1) + I(Y; \hat{Y}) \\ I(Y; T_{\text{MLP}}) &= I(Y; \hat{Y}) \end{aligned} \quad (15)$$

209 where T_{MLP} denotes a random variable corresponding to the entire architecture of the MLP.

210 *Proof:* Since $H(Y) = I(X; Y)$ (Appendix D), $H(X) = H(\bar{X}) + I(X; Y) = H(\bar{X}) + H(Y)$.

211 Hence, Corollary 2 can be derived by Corollary 1 and the chain rule. The proof is in Appendix F.

212 4 Simulations

213 In this section, we propose a mutual information estimator based on $(\Omega_T, \mathcal{F}, P_T)$ and demonstrate the
 214 probabilistic representation and information theoretic explanations for deep learning on a synthetic
 215 dataset with known entropy. Additional experiments on benchmark datasets are in Appendix H.

216 4.1 Setup

217 **Mutual information estimator.** Based on the definition of mutual information, we have

$$I(X; T_i) = H(T_i) - H(T_i|X). \quad (16)$$

218 Previous works simply estimate $I(X; T_i) = H(T_i)$, because T_i is assumed to be entirely dependent
 219 on X in the Markov chain, Equation (2), thereby $H(T_i|X) = 0$. However, Corollary 1 shows that T_i
 220 depends on both X and Y if taking into account the training phase, thereby $H(T_i|X) \neq 0$.

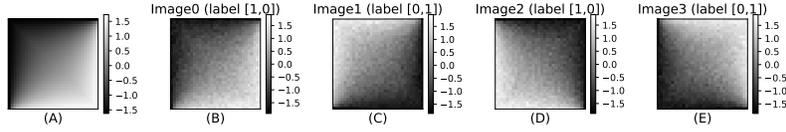


Figure 3: (A) the deterministic image $\hat{\mathbf{x}}$. Image0 is generated by adding $\mathcal{N}(\mu, \sigma^2)$ without rotation, Image1 is generated by rotating $\hat{\mathbf{x}}$ along the secondary diagonal direction and adding $\mathcal{N}(\mu, \sigma^2)$, Image2 and Image are generated by rotating $\hat{\mathbf{x}}$ along the vertical and horizontal directions, respectively, and adding $\mathcal{N}(\mu, \sigma^2)$.

Table 1: The number of neurons(nodes) and the activation function in the layers of the MLPs

	\mathbf{x}	\mathbf{t}_1	\mathbf{t}_2	$\hat{\mathbf{y}}$	$\sigma(\cdot)$
MLP1	1024 (32×32)	8	6	2	$\text{ReLU}(z) = \max(0, z)$
MLP2	1024 (32×32)	8	6	2	$\text{Tanh}(z) = (e^z - e^{-z}) / (e^z + e^{-z})$
MLP3	1024 (32×32)	2	6	2	ReLU

221 To accurately estimate $I(X; T_i)$, we need to specify $P(T_i|X)$ and $P(T_i)$. Based on $(\Omega_{T_i}, \mathcal{F}, P_{T_i})$,
 222 we formulate $P_{T_i|X}(n|\mathbf{x}^j)$ of the three fully connected layers in the MLP as

$$\begin{aligned}
 P_{T_1|X}(n|\mathbf{x}^j) &= \frac{1}{Z_{F_1}} \exp[\sigma_1(\langle \omega_n^{(1)}, \mathbf{x}^j \rangle)], \quad P_{T_2|X}(k|\mathbf{x}^j) = \frac{1}{Z_{F_2}} \exp[\sigma_2(\langle \omega_k^{(2)}, \mathbf{t}_1(\mathbf{x}^j) \rangle)], \\
 P_{T_Y|X}(l|\mathbf{x}^j) &= \frac{1}{Z_{F_Y}} \exp[\langle \omega_l^{(3)}, \mathbf{t}_2(\mathbf{t}_1(\mathbf{x}^j)) \rangle].
 \end{aligned}
 \tag{17}$$

223 To derive the marginal distribution $P(T_i)$, we sum the joint distribution $P(T_i, X)$ over $\mathbf{x} \in \mathcal{X}$,

$$P(T_i = n) = \sum_{\mathbf{x} \in \mathcal{X}} P_X(\mathbf{x}) P_{T_i|X}(n|\mathbf{x}) \approx \sum_{\mathbf{x}^j \in \mathcal{D}} P_X(\mathbf{x}^j) P_{T_i|X}(n|\mathbf{x}^j) = \frac{1}{J} \sum_{\mathbf{x}^j \in \mathcal{D}} P_{T_i|X}(n|\mathbf{x}^j),
 \tag{18}$$

224 where $P_X(\mathbf{x}^j)$ is estimated by the empirical distribution $1/J$ given \mathcal{D} . Finally, we can derive $I(X; T_i)$
 225 by Equation 16, 17, and 18. Similarly, based on the definition of mutual information, we have

$$I(Y; T_i) = H(T_i) - H(T_i|Y).
 \tag{19}$$

226 To estimate $H(T_i|Y)$, we reformulate $P(T_i|Y)$ as

$$P_{T_i|Y}(n|l) = \sum_{\mathbf{x} \in \mathcal{X}} P_{T_i|X}(n|\mathbf{x}) P_{X|Y}(\mathbf{x}|l) \approx \frac{1}{N(l)} \sum_{\mathbf{x}^j \in \mathcal{D}, y^j=l} P_{T_i|X}(n|\mathbf{x}^j),
 \tag{20}$$

227 where $P_{X|Y}(\mathbf{x}^j|l)$ is estimated by the empirical distribution $1/N(l)$ and $N(l)$ denotes the number of
 228 samples with the label l in \mathcal{D} . Finally, we can derive $I(Y; T_i)$ by Equation 18, 19, and 20.

229 **Synthetic dataset.** The dataset consists of 512 gray-scale 32×32 images, which are evenly generated
 230 by rotating a deterministic image $\hat{\mathbf{x}}$ in four different orientations and adding Gaussian noise with
 231 expectation $\mu = \mathbb{E}(\hat{\mathbf{x}})$ and variance $\sigma^2 = 1$, namely $\mathbf{x} = r(\hat{\mathbf{x}}) + \mathcal{N}(\mu, \sigma^2)$, where $r(\cdot)$ denotes
 232 the rotation method shown in Figure 3. The reason for adding Gaussian noise is to avoid DNNs
 233 directly memorizing the deterministic image. In addition, the binary labels [1,0] and [0,1] evenly
 234 divide the synthetic dataset into two classes. As a result, the synthetic dataset has (approximately)
 235 2 bits information and the labels have 1 bit information. Compared to popular benchmark dataset
 236 with unknown features and entropy, *e.g.*, MNIST [17] and Fashion-MNIST [36], the features and the
 237 entropy of the synthetic dataset are clear and known, which enables us to examine the probabilistic
 238 representation and the mutual information estimator.

239 **Neural Networks.** We train three MLPs, namely MLP1, MLP2 and MLP3, on the synthetic dataset
 240 by a variant of Stochastic Gradient Descent (SGD) method, namely Adam [13], over 1000 epochs
 241 with the learning rate $\alpha = 0.03$. Table 1 summarizes the architecture of the three MLPs.

242 4.2 Validating the probability space and the mutual information estimator

243 We demonstrate the sample space Ω_T by visualizing the weights⁴ of the eight neurons in \mathbf{t}_1 , *i.e.*,
 244 $\omega_n^{(1)} = \{\omega_{mn}^{(1)}\}_{m=1}^{1024}$, in 5 different epochs (*i.e.*, 0,1,4,128,1000) in Figure 4 (Left). As training
 245 continues, we observe that $\omega_n^{(1)}$ quickly learns all the spatial features of the synthetic dataset. For
 246 instance, $\omega_2^{(1)}$ has low magnitude at top-left positions and high magnitude at bottom-right positions,
 247 which correctly characterizes the spatial feature of Image0. Similarly, $\omega_3^{(1)}$, $\omega_4^{(1)}$, and $\omega_5^{(1)}$ correctly
 248 characterize the spatial feature of Image1, Image2, and Image3 in Figure 3, respectively.

⁴We only show the learned weights in MLP1 because we observe that the learned weights in MLP1 and MLP2 are very similar, though they use different activation functions.

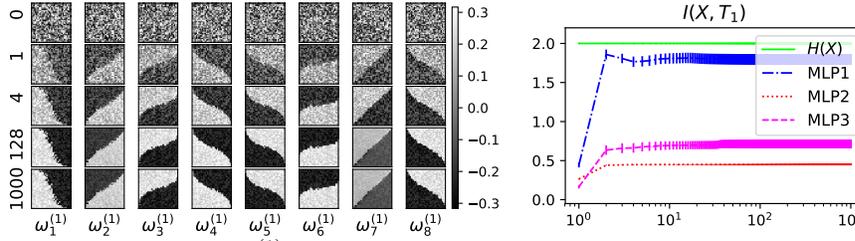


Figure 4: (Left) The eight features $\{\omega_n^{(1)}\}_{n=1}^8$ learned by the weights of the eight neurons in 5 different epochs (i.e., 0, 1, 4, 128, 1000), where $\omega_n^{(1)} = \{\omega_{mn}^{(1)}\}_{m=1}^{1024}$ are reshaped into 32×32 to show the spatial structure. (Right) The variation of $I(X; T_1)$ in the MLP1, MLP2, and MLP3 during 1000 epochs.

Table 2: The Gibbs probability $P_{F_1|X}(\omega_n^{(1)}|\text{Image0})$ in MLP1 and MLP2 in the 1000 epoch

	$\omega_1^{(1)}$	$\omega_2^{(1)}$	$\omega_3^{(1)}$	$\omega_4^{(1)}$	$\omega_5^{(1)}$	$\omega_6^{(1)}$	$\omega_7^{(1)}$	$\omega_8^{(1)}$
$\langle \omega_n^{(1)}, \mathbf{x} \rangle$	-63.6	208.8	-181.6	45.1	-55.6	157.5	-210.0	-30.1
$f_{1n}^{\text{ReLU}}(\mathbf{x})$	0.0	208.8	0.0	45.1	0.0	157.5	0.0	0.0
$\exp[f_{1n}^{\text{ReLU}}(\mathbf{x})]$	1.0	4.79e+90	1.0	3.86e+19	1.0	2.51e+68	1.0	1.0
$P_{T_1 X}^{\text{ReLU}}$	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0
$f_{1n}^{\text{Tanh}}(\mathbf{x})$	-1.0	1.0	-1.0	1.0	-1.0	1.0	-1.0	-1.0
$\exp[f_{1n}^{\text{Tanh}}(\mathbf{x})]$	0.36	2.71	0.36	2.71	0.36	2.71	0.36	0.36
$P_{T_1 X}^{\text{Tanh}}$	0.037	0.272	0.037	0.272	0.037	0.272	0.037	0.037

$f_{1n}^{\text{Tanh}}(\mathbf{x}) = \sigma^{\text{Tanh}}(\langle \omega_n^{(1)}, \mathbf{x} \rangle)$ and $f_{1n}^{\text{ReLU}}(\mathbf{x}) = \sigma^{\text{ReLU}}(\langle \omega_n^{(1)}, \mathbf{x} \rangle)$ are the activations given the same $\langle \omega_n^{(1)}, \mathbf{x} \rangle$.

249 We demonstrate that $P(T_1|X)$ correctly measures the probability of $\{\omega_n^{(1)}\}_{n=1}^8$ being recognized the
 250 feature with the largest cross-correlation to \mathbf{x} in Table 2. For instance, $\omega_2^{(1)}$ correctly characterizes
 251 the feature of Image0 and has the largest cross-correlation $\langle \omega_2^{(1)}, \mathbf{x} \rangle = 190.8$, thus it has the largest
 252 probability $P_{T_1|X}^{\text{ReLU}}(\omega_2^{(1)}|\text{Image0}) = 1.0$ being recognized as the feature with largest cross-correlation
 253 to Image0. In contrast, since $\omega_7^{(1)}$ incorrectly characterizes the feature of Image0 and has the lowest
 254 cross-correlation $\langle \omega_7^{(1)}, \mathbf{x} \rangle = -210.0$, so it has the lowest probability $P_{T_1|X}^{\text{ReLU}}(\omega_7^{(1)}|\text{Image0}) = 0.0$
 255 being recognized as the feature with largest cross-correlation to Image0.

256 We observe that an activation function (abbr. ACT) plays an important role in the distribution.
 257 Specifically, ReLU, a non-saturating (unbounded) ACT [9], preserves the positive cross-correlations
 258 while resets all the negative ones as zero. $P_{T_1|X}^{\text{ReLU}}(\omega_2^{(1)}|\text{Image0}) = 1.0$ shows that ReLU derives the
 259 correct probability of $\omega_2^{(1)}$ being recognized as the feature with largest cross-correlation. In contrast,
 260 though $\omega_2^{(1)}$ has stronger cross-correlation to Image0 than $\omega_4^{(1)}$, i.e., $\langle \omega_2^{(1)}, \mathbf{x} \rangle > \langle \omega_4^{(1)}, \mathbf{x} \rangle$, Tanh, a
 261 saturating (bounded) ACT, derives $f_{12}^{\text{Tanh}}(\mathbf{x}) = f_{14}^{\text{Tanh}}(\mathbf{x}) = 1.0$, and makes $\omega_4^{(1)}$ to incorrectly have
 262 the same probability 0.272 to $\omega_2^{(1)}$ being recognized as the feature with the largest cross-correlation
 263 to Image0, i.e., Tanh hinders t_1 from correctly recognizing the features of input. The simulations for
 264 validating the probability space based on other synthetic images are presented in Appendix G.

265 To validate the mutual information estimator, we follow recent works [30, 28] to train the three
 266 MLPs with 50 different random initialization and study the average mutual information. Figure 4
 267 (Right) shows that $I(X; T_1)$ quickly increases to 1.81 and keeps stable in the MLP1, i.e., t_1 learns
 268 most information of the dataset as $H(X) = 2.0$. Notably, the result is consistent with the variation
 269 of the weights in Figure 4 (Left), which shows that the weights correctly characterize the features
 270 of the dataset and keeps stable after the fourth epoch. As a comparison, we observe that $I(X; T_1)$
 271 keeps stable at 0.44 in the MLP2, which confirms the statement that Tanh hinders t_1 from correctly
 272 recognizing the features of input. In addition, Figure 4 (Right) shows that $I(X; T_1) \approx 0.79$ in MLP3
 273 is smaller than $I(X; T_1) \approx 1.81$ in MLP1, which is consistent with Definition 1, i.e., a layer with
 274 fewer neurons would represent fewer possible features, thus it contains less information.

275 In summary, we demonstrate the probability space $(\Omega_T, \mathcal{F}, P_T)$ and show that if an ACT cannot
 276 preserve the cross-correlation between weights(features) and input, it would distort the distribution
 277 of a layer, thereby affecting the mutual information between the layer and data/labels. In addition,
 278 we show that the proposed mutual information estimator outperforms the existing non-parametric
 279 models, e.g., empirical distribution [30] and KDE [28], based on the synthetic dataset. Especially,
 280 activations do not satisfy the *i.i.d.* prerequisite of non-parametric models is an important reason for
 281 non-parametric models deriving inaccurate mutual information in DNNs. Due to limited space, the
 282 experimental comparison and study of non-parametric models are presented in Appendix G.

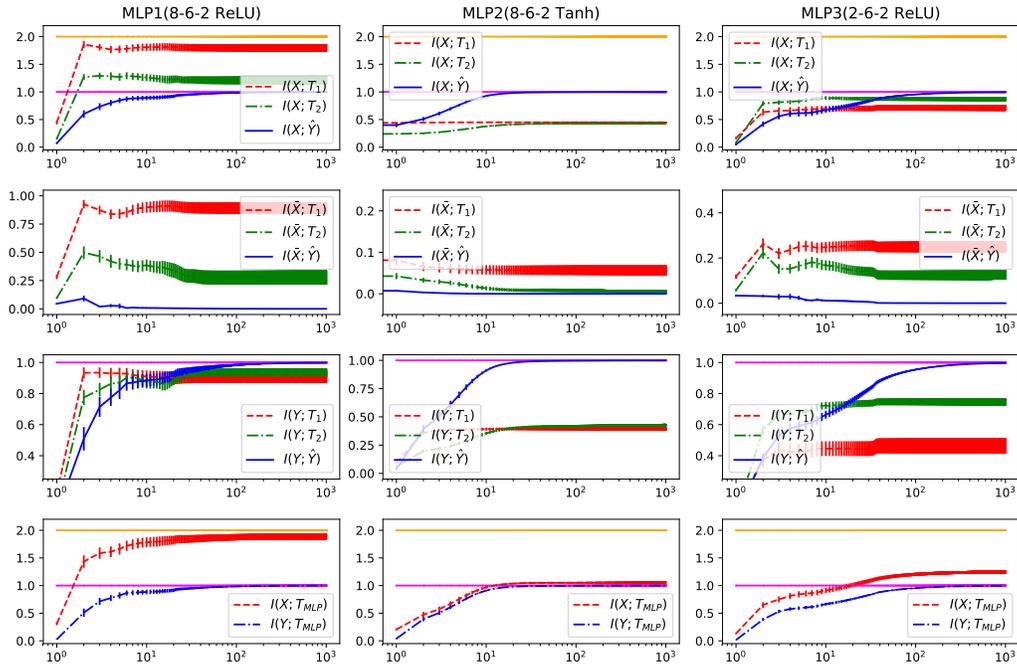


Figure 5: All the x-axis index training epochs. In each column, the first three figures show $I(X; T_i)$, $I(\bar{X}; T_i)$, and $I(Y; T_i)$ respectively. The fourth figure shows $I(X; T_{MLP})$ and $I(Y; T_{MLP})$ in a MLP. The pink line denotes $H(Y) = 1.0$ and the orange line denotes $H(X) = 2.0$.

283 4.3 Validating the information theoretic explanations for DNNs

284 In Figure 5, we observe $I(X; T_i) \leq I(X; \hat{Y})$ in MLP2 and MLP3, which confirms that the Markov
 285 chain proposed by previous works, Equation (2), cannot fully explain the information flow in MLPs,
 286 if taking into account the back-propagation training. As a comparison, the second and third row
 287 show $I(\bar{X}; T_1) \geq I(\bar{X}; T_2) \geq I(\bar{X}; \hat{Y})$ and $I(Y; T_1) \leq I(Y; T_2) \geq I(Y; \hat{Y})$ in all the three MLPs,
 288 which validates that Corollary 1, *i.e.*, Equation (14) characterizes the information flow in MLPs.

289 Figure 5 demonstrates that different hidden layers achieve different IB trade-offs depending on
 290 the architecture and the position of the layers in MLPs. In terms of architecture, $I(Y; T_1) > 0.8$
 291 and $I(\bar{X}; T_1) > 0.75$ in MLP1 indicate that t_1 , with ReLU, achieves a good prediction without
 292 much compression, whereas $I(Y; T_1) < 0.5$ and $I(\bar{X}; T_1) < 0.1$ in MLP2 show that t_1 , with Tanh,
 293 achieves a different IB trade-off. In addition, $I(Y; T_1) \approx 0.45$ and $I(\bar{X}; T_1) \approx 0.25$ in MLP3 show
 294 the effect of neuron numbers on the IB trade-off. In terms of position, $I(Y; \hat{Y}) = 1$ and $I(\bar{X}; \hat{Y}) = 0$
 295 in MLP1 means that \hat{y} has a different IB trade-off to t_1 in MLP1.

296 We demonstrate that a MLP satisfies the IB principle no matter what the architecture of the MLP
 297 is. Figure 5 visualizes $I(X; T_{MLP})$ and $I(Y; T_{MLP})$ based on Corollary 2. It shows that all of three
 298 MLPs satisfy the IB principle, namely $I(X; T_{MLP}) < H(X) = 2$ and $I(Y; T_{MLP}) = H(Y) = 1$,
 299 though they have different architectures. Importantly, in contrast to previous work [28] claiming that
 300 the compression not exists in DNNs with non-saturating ACT, such as ReLU, Figure 5 clearly shows
 301 that the compression exists in all the MLPs, no matter the activation function of MLPs.

302 We further demonstrate the information theoretic explanations for DNNs on the benchmark MNIST
 303 and Fashion-MNIST datasets. The experiments are presented in Appendix H.

304 5 Conclusion and future work

305 In this work, we (1) specify the probability space for a hidden layer for (2) accurately estimating the
 306 mutual information and (3) clearly explaining how the components of the layer affect the mutual
 307 information. We take into account the back-propagation training and derive two novel Markov chains
 308 to characterize the information flow in DNNs. Furthermore, we demonstrate that a DNN satisfies the
 309 IB principle no matter the architecture of the DNN. In contrast, different hidden layers show different
 310 IB trade-offs depending on the architecture and the position of the layers in DNNs. A potential
 311 direction is to study the generalization of DNNs based on the probabilistic representation.

References

- 312
- 313 [1] Rana Ali Amjad and Bernhard Claus Geiger. Learning representations for neural network-based classi-
314 fication using the information bottleneck principle. *IEEE transactions on pattern analysis and machine*
315 *intelligence*, 2019.
- 316 [2] George Casella and Roger L Berger. *Statistical inference*. Cengage Learning, 2021.
- 317 [3] Ivan Chelombiev, Conor Houghton, and Cian O’Donnell. Adaptive estimators show information compres-
318 sion in deep neural networks. In *International Conference on Learning Representations*, 2019.
- 319 [4] Thomas Cover and Joy Thomas. *Elements of Information Theory*. Wiley-Interscience, Hoboken, New
320 Jersey, 2006.
- 321 [5] Balázs Csanád Csáji et al. Approximation with artificial neural networks. *Faculty of Sciences, Eötvös Lornd*
322 *University, Hungary*, 24(48):7, 2001.
- 323 [6] Rick Durrett. *Probability: theory and examples*, volume 49. Cambridge university press, 2019.
- 324 [7] Marylou Gabrié, Andre Manoel, Clément Luneau, Nicolas Macris, Florent Krzakala, Lenka Zdeborová,
325 et al. Entropy and mutual information in models of deep neural networks. In *Advances in Neural*
326 *Information Processing Systems*, pages 1821–1831, 2018.
- 327 [8] Matt W Gardner and SR Dorling. Artificial neural networks (the multilayer perceptron)—a review of
328 applications in the atmospheric sciences. *Atmospheric environment*, 32(14-15):2627–2636, 1998.
- 329 [9] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural
330 networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*,
331 pages 249–256, 2010.
- 332 [10] Ziv Goldfeld, Ewout Van Den Berg, Kristjan Greenewald, Igor Melnyk, Nam Nguyen, Brian Kingsbury,
333 and Yury Polyanskiy. Estimating information flow in deep neural networks. In *Proceedings of the 36th*
334 *International Conference on Machine Learning*, volume 97, pages 2299–2308, 2019.
- 335 [11] Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT
336 press Cambridge, 2016.
- 337 [12] Patrick Kidger and Terry Lyons. Universal approximation with deep narrow networks. In *Conference on*
338 *Learning Theory*, pages 2306–2327. PMLR, 2020.
- 339 [13] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
340 *arXiv:1412.6980*, 2014.
- 341 [14] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In
342 *International Conference on Machine Learning*, pages 1885–1894. PMLR, 2017.
- 343 [15] Artemy Kolchinsky, Brendan D Tracey, and Steven Van Kuyk. Caveats for information bottleneck in
344 deterministic scenarios. *arXiv preprint arXiv:1808.07593*, 2018.
- 345 [16] Alex Krizhevsky. Learning multiple layers of features from tiny images. *Technique Report*, 2009.
- 346 [17] Y. LeCun, BE. Boser, and JS. Denker. Handwritten digit recognition with a back-propagation network. In
347 *NeurIPS*, pages 396–494, 1990.
- 348 [18] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
- 349 [19] Yann LeCun, Sumit Chopra, Raia Hadsell, Marc’Aurelio Ranzato, and Fu Jie Huang. *A tutorial on*
350 *energy-based learning*. MIT Press, 2006.
- 351 [20] Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S. Schoenholz, Jeffrey Pennington, and Jascha
352 Sohl-Dickstein. Deep neural networks as gaussian processes. In *ICLR*, 2018.
- 353 [21] Aravindh Mahendran and Andrea Vedaldi. Understanding deep image representations by inverting them.
354 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5188–5196,
355 2015.
- 356 [22] David McAllester and Karl Stratos. Formal limitations on the measurement of mutual information. In
357 *International Conference on Artificial Intelligence and Statistics*, pages 875–884. PMLR, 2020.

- 358 [23] Julian D Olden and Donald A Jackson. Illuminating the “black box”: a randomization approach for
359 understanding variable contributions in artificial neural networks. *Ecological modelling*, 154(1-2):135–150,
360 2002.
- 361 [24] Liam Paninski. Estimation of entropy and mutual information. *Neural computation*, 15(6):1191–1253,
362 2003.
- 363 [25] David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. Learning representations by back-
364 propagating errors. *Nature*, 323:533–536, October 1986.
- 365 [26] Ruslan Salakhutdinov and Geoffrey Hinton. Deep boltzmann machines. In *Artificial intelligence and*
366 *statistics*, pages 448–455. PMLR, 2009.
- 367 [27] Wojciech Samek, Thomas Wiegand, and Klaus-Robert Müller. Explainable artificial intelligence: Under-
368 standing, visualizing and interpreting deep learning models. *arXiv preprint arXiv:1708.08296*, 2017.
- 369 [28] Andrew Saxe, Yamini Bansal, Joel Dapello, Madhu Advani, Artemy Kolchinsky, Brendan Tracey, and
370 David Cox. On the information bottleneck theory of deep learning. In *International Conference on*
371 *Representation Learning*, 2018.
- 372 [29] Ohad Shamir, Sivan Sabato, and Naftali Tishby. Learning and generalization with the information
373 bottleneck. *Theoretical Computer Science*, 411(29-30):2696–2711, 2010.
- 374 [30] Ravid Shwartz-Ziv and Naftali Tishby. Opening the black box of deep neural networks via information.
375 *arXiv preprint arXiv:1703.00810*, 2017.
- 376 [31] Noam Slonim. *The information bottleneck: Theory and applications*. PhD thesis, Citeseer, 2002.
- 377 [32] Naftali Tishby, Fernando C Pereira, and William Bialek. The information bottleneck method. *arXiv*
378 *preprint physics/0004057*, 2000.
- 379 [33] Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. In *2015 IEEE*
380 *Information Theory Workshop (ITW)*, pages 1–5. IEEE, 2015.
- 381 [34] Vladimir N Vapnik. An overview of statistical learning theory. *IEEE transactions on neural networks*,
382 10(5):988–999, 1999.
- 383 [35] Larry Wasserman. *All of nonparametric statistics*. Springer Science & Business Media, 2006.
- 384 [36] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking
385 machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 386 [37] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European*
387 *conference on computer vision*, pages 818–833. Springer, 2014.
- 388 [38] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep
389 learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.

390 **Checklist**

- 391 1. For all authors...
- 392 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contribu-
393 tions and scope? [Yes]
- 394 (b) Did you describe the limitations of your work? [Yes] see Section 5
- 395 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 396 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 397 2. If you are including theoretical results...
- 398 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 399 (b) Did you include complete proofs of all theoretical results? [Yes]
- 400 3. If you ran experiments...
- 401 (a) Did you include the code, data, and instructions needed to reproduce the main experimental
402 results (either in the supplemental material or as a URL)? [Yes] see the URL in Appendix G
- 403 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)?
404 [Yes] see Section 4.1, Appendix G, and Appendix H
- 405 (c) Did you report error bars (e.g., with respect to the random seed after running experiments
406 multiple times)? [Yes] see Section 4
- 407 (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs,
408 internal cluster, or cloud provider)? [Yes] see Appendix G
- 409 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 410 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 411 (b) Did you mention the license of the assets? [Yes]
- 412 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] see
413 Appendix H
- 414 (d) Did you discuss whether and how consent was obtained from people whose data you’re us-
415 ing/curating? [N/A]
- 416 (e) Did you discuss whether the data you are using/curating contains personally identifiable informa-
417 tion or offensive content? [N/A]
- 418 5. If you used crowdsourcing or conducted research with human subjects...
- 419 (a) Did you include the full text of instructions given to participants and screenshots, if applicable?
420 [N/A]
- 421 (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB)
422 approvals, if applicable? [N/A]
- 423 (c) Did you include the estimated hourly wage paid to participants and the total amount spent on
424 participant compensation? [N/A]